

Question #1Monte-Carlo Policy Iteration : (Every-Visit)Pseudocode:

$$Q(s, a) = 0, \pi(s) = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}, \gamma = 0.9$$

Episode-length = 5 , Returns(s,a) = [] $\forall s, a$

Repeat until (policy unchanged from last iteration)

- Generate an episode using π
- For ~~each~~^{every} pair (s, a) in episode :
- $R \leftarrow$ return following first occurrence of s, a
- Append R to Returns(s, a)
- $Q(s, a) \leftarrow \text{Avege}(\text>Returns(s, a))$

- For each s in the episode

- $\pi(s) \leftarrow \underset{a}{\text{argmax}} Q(s, a)$

let

- When generating episodes, let the starting state be fixed at A. For exploration, let the agent follow an ϵ -greedy policy with $\epsilon = 0.9$.

Pseudocode for generating an episode:

procedure $\text{gen-episode}(\text{in Ep-length}, \pi, E$
out episode)

- $s = A$, episode = []
- while len(episode) < Ep-length
- $a = \begin{cases} \pi(s) & 1 - E \\ a' \text{ or } a'' \text{ randomly } E \end{cases}$
- take a, observe r, move to s'
- append (s, a, s', r) to episode
- $s = s'$
- return episode.

Iteration 1:

- episode = gen-episode(...)
 $\hookrightarrow [(0,0,0,0), (0,1,1,4), (1,1,0,-1),$
 $(0,0,0,0), (0,0,0,0)]$

My mapping is : $A \rightarrow 0$

$B \rightarrow 1$

$\underbrace{}$

states

$a' \rightarrow 0$

$a'' \rightarrow 1$

$\underbrace{}$

actions

which means:

\rightarrow episode'	s	a	s'	r
1	A	a'	A	0
2	A	a^2	B	4
3	B	a^2	A	-1
4	A	a'	A	0
5	A	a'	A	0

$\rightarrow (s, a)$ pairs = (A, a') , (A, a^2) , (B, a^2) ,

1) for (A, a') :

$$\rightarrow R = 0 + 0.9 \times 4 + 0.9^2 \times -1 + 0 + 0$$

$$\rightarrow R = 2.79$$

$$\rightarrow \text{Returns}(s, a) = \text{Returns}(A, a') = [2.79]$$

$$\rightarrow Q(s, a) = Q(A, a') = 2.79$$

2) for (A, a^2) :

$$\rightarrow R = 4 + 0.9 \times -1 = 3.1$$

$$\rightarrow \text{Returns}(A, a^2) = [3.1]$$

$$\rightarrow Q(A, a^2) = 3.1$$

3) for (B, a^2) :

$$\rightarrow R = -1$$

$$\rightarrow \text{Returns}(B, a^2) = [-1]$$

$$\rightarrow Q(B, a^2) = -1$$

4) for (A, a') :

$$\rightarrow R = 0$$

$$\rightarrow \text{Returns}(A, a') = [2.79, 0]$$

$$\rightarrow Q(sA, a') = (2.79 + 0) / 2 = 1.395$$

5) for (A, a') :

$$\rightarrow R = 0$$

$$\rightarrow \text{Returns}(A, a') = [2.79, 0, 0]$$

$$\rightarrow Q(A, a') = (2.79 + 0 + 0) / 3$$

$$\Rightarrow Q(A, a') = 0.93$$

All (s, a) pairs visited.

$$Q = \begin{matrix} & A & [0.93 & 3.1] \\ & B & [0 & -1] \\ a' & & a^2 \end{matrix}$$

$$\text{Returns} = \begin{matrix} & A & [[2.79, 0, 0] & [3.1]] \\ & B & [[] & [-1]] \\ a' & & a^2 \end{matrix}$$

$$\pi(s=A) = \underset{a}{\text{argmax}} (0.93, 3.1) = a^2$$

$$\pi(s=B) = \underset{a}{\text{argmax}} (0, -1) = a'$$

$$\text{Thus, } \pi' = \begin{bmatrix} a^2 \\ a' \end{bmatrix}$$

Date: _____

Iteration 2:

→ episode = gen-episode (...)
 → episode²

	s	a	s'	r
1	A	a ²	B	4
2	B	a'	B	5
3	B	a'	B	5
4	B	a'	B	5
5	B	a'	B	5

1) for (A, a²):

$$\rightarrow R = 4 + 0.9 \times 5 + 0.9^2 \times 5 + 0.9^3 \times 5 + 0.9^4 \times 5 \\ R = 19.4755$$

$$\rightarrow \text{Returns}(A, a^2) = [3.1, 19.4755]$$

$$\rightarrow Q(A, a^2) = (3.1 + 19.4755) / 2$$

$$Q(A, a^2) = 11.28755$$

2) for (B, a¹):

$$\rightarrow R = 5 + 0.9 \times 5 + 0.9^2 \times 5 + 0.9^3 \times 5 \\ R = 17.195$$

$$\rightarrow \text{Returns}(B, a^1) = [17.195]$$

$$\rightarrow Q(B, a^1) = 17.195$$

3) for (B, a¹):

$$\rightarrow R = 5 + 0.9 \times 5 + 0.9^2 \times 5$$

$$R = 13.55$$

$$\rightarrow \text{Returns}(B, a^1) = [17.195, 13.55]$$

$$\rightarrow Q(B, a^1) = (17.195 + 13.55) / 2$$

$$Q(B, a') = 15.3725$$

4) for (B, a') :

$$\rightarrow R = 5 + 0.9 \times 5 = 9.5$$

$$\rightarrow \text{Returns}(B, a') = [17.95, 13.55, 9.5]$$

$$\rightarrow Q(B, a') = 13.915$$

5) for (B, a') :

$$\rightarrow R = 5$$

$$\rightarrow \text{Returns}(B, a') = [17.95, 13.55, 9.5, 5]$$

$$\rightarrow Q(B, a') = 11.31125$$

All (s, a) pairs visited.

$$Q = A \begin{bmatrix} 0.93 & 11.28775 \\ B & \end{bmatrix}$$

$$\begin{array}{c} a' \\ a'' \end{array}$$

$$\text{Returns} = \begin{bmatrix} [2.79, 0, 0] & [3.1, 19.4755] \\ [17.95, 13.55, 9.5, 5] & [-1] \end{bmatrix}$$

$$\pi(s=A) \Rightarrow \underset{a}{\text{argmax}} (0.93, 11.28) = a^*$$

$$\begin{array}{c} a' \\ a'' \end{array}$$

$$\pi(s=B) = \underset{a}{\text{argmax}} (11.31, -1) = a'$$

$$\begin{array}{c} a' \\ a'' \end{array}$$

$$\text{Thus, } \pi = \begin{bmatrix} a'' & A \\ a' & B \end{bmatrix}$$

Policy π^* is the same as π' , so we terminate.

Question 2

Tabular Actor-Critic

Pseudocode:

$V(s) = 0, H(s, a) = 0$ for all $s \in S, a \in A$

Repeat for N episodes:

$t = 0, s_0 \in S$ (random)

while $t <$ episode-length:

$$\cdot \pi(a|s) = \frac{e^{H(s,a)}}{\sum_{a \in A} e^{H(s,a)}} \quad \forall s, a$$

- select action: $a_t \sim \pi(\cdot | s_t)$
- take action a_t , move to s_{t+1} , observe r_{t+1}
- $\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$
- $V(s_t) = V(s_t) + \alpha \delta_t$
- $H(s_t, a_t) = H(s_t, a_t) + \beta \delta_t (1 - \pi(a_t | s_t))$
- $t = t + 1$

We are provided with an episode's state-action-reward information, and require computing intermediate policies (π), state values (V), and preferences (H).

Date: _____

$$\alpha = 0.5, \beta = 0.1, r = 0.9$$

$$V(s) = \begin{bmatrix} 0 & 0 \\ A & B \end{bmatrix}$$

$$H(s, a) = \begin{array}{|c|cc|} \hline & 0 & 0 \\ \hline A & 0 & 0 \\ \hline a' & a^* & \\ \hline \end{array}$$

$$t = 0:$$

$$\rightarrow s_t = A, a_t = a^*, r_{t+1} = 10$$

$$s_{t+1} = A$$

$$\pi(a^* | s=A)$$

$$\rightarrow \pi(a=a^* | s=A) = \frac{e^{H(A, a^*)}}{e^{H(A, a^*)} + e^{H(A, a')}} = \frac{e^0}{e^0 + e^0}$$

$$\pi(a^* | A) = 0.5$$

Similarly,

$$\pi(a^* | B) = \pi(a^* | A) = \pi(a^* | B) = 0.5$$

$$\rightarrow S_t = 10 + 0.9 \times V(A) - V(A) = 10 + 0 - 0$$

$$S_t = 10$$

$$\rightarrow V(A) = V(A) + 0.5 \times 10 = 5$$

$$\rightarrow H(A, a^*) = 0 + 0.1 \times 10 \times (1 - \pi(a^* | A)) \\ = 0.1 \times 10 \times (1 - 0.5)$$

$$\rightarrow H(A, a^*) = 0.5$$

Date: _____

ω_0 ,

$$\pi^0 = A \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$V^0 = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}$$

$$H^0 = A \begin{bmatrix} 0.5 & 0 \\ 0 & 0 \end{bmatrix}$$

$$t=1: \quad \begin{matrix} a' & a' \end{matrix}$$

$$\rightarrow S_t = A, a_t = a^2, r_{t+1} = -5, S_{t+1} = B$$

$$\rightarrow \pi(a' | A) = e^{H(A, a')} / (e^{H(A, a')} + e^{H(A, a^2)})$$

$$= e^{0.5} / (e^{0.5} + e^0)$$

$$\Rightarrow \boxed{\pi(a' | A) = 0.622}$$

Similarly,

$$\pi(a^2 | A) = e^0 / (e^{0.5} + e^0)$$

$$\boxed{\pi(a^2 | A) = 0.377}$$

$\pi(\cdot | B)$ remain the same.

$$\rightarrow \delta_t = -5 + 0.9 \times V(B) - V(A)$$

$$S_t = -5 + 0.9 \times 0 - 5 = -10$$

$$\rightarrow V(A) = 5 + 0.5 \times -10 = 5 - 5$$

$$\boxed{V(A) = 0}$$

$$\rightarrow H(A, a^2) = 0 + 0.1 \times -10 \times (1 - \pi(a^2 | A))$$

$$= 0.1 \times -10 \times (1 - 0.377)$$

$$\boxed{H(A, a^2) = -0.623}$$

Date: _____

So,

$$\pi' = \begin{bmatrix} 0.622 & 0.377 \\ 0.5 & 0.5 \end{bmatrix}$$

$$V' = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$H' = \begin{bmatrix} 0.5 & -0.623 \\ 0 & 0 \end{bmatrix}$$

$t=2:$

$$\rightarrow a_1 = B, a_2 = a', \delta_{t+1} = 40, S_{t+1} = A$$

$$\rightarrow \pi(a_1 | A) = e^{0.5} / (e^{0.5} + e^{-0.623})$$

$$\boxed{\pi(a_1 | A) = 0.7545}$$

$$\pi(a_2 | A) = e^{-0.623} / (e^{0.5} + e^{-0.623})$$

$$\boxed{\pi(a_2 | A) = 0.245}$$

$\pi(\cdot | B)$ unchanged.

$$\rightarrow S_t = 40 + 0.9 \times V(A) - V(B)$$

$$S_t = 40 + 0.9 \times 0 - 0$$

$$S_t = 40$$

$$\rightarrow V(B) = 0 + 0.5 \times 40$$

$$\boxed{V(B) = 20}$$

$$\rightarrow H(B, a') = 0 + 0.1 \times 40 \times (1 - 0.5)$$

$$\boxed{H(B, a') = 2}$$

So,

$$\pi = \begin{bmatrix} 0.7545 & 0.245 \\ 0.5 & 0.5 \end{bmatrix}$$

$$V = \begin{bmatrix} 0 & 20 \end{bmatrix}$$

$$H = \begin{bmatrix} 0.5 & -0.623 \\ 2 & 0 \end{bmatrix}$$

Date:

$t=3$:

$$\rightarrow s_t = A, a_t = a^2, \gamma_{t+1} = -5, s_{t+1} = B$$

$$\rightarrow \pi(a^1 | B) = e^2 / (e^2 + e^0)$$

$$\boxed{\pi(a^1 | B) = 0.88}$$

$$\boxed{\pi(a^2 | B) = e^0 / (e^0 + e^2)}$$

$$\boxed{\pi(a^2 | B) = 0.12}$$

$\pi(\cdot | A)$ unchanged

$$\rightarrow \delta_t = -5 + 0.9 \times V(B) - V(A)$$

$$= -5 + 0.9 \times 20 - 0 = 13$$

$$\rightarrow V(A) = 0 + 0.5 \times 13$$

$$\boxed{V(A) = 6.5}$$

$$\rightarrow H(A, a^2) = -0.623 + 0.1 \times 13 \times (1 - 0.245)$$

$$\boxed{H(A, a^2) = 0.358}$$

so,

$$\pi^3 = \begin{bmatrix} 0.7545 & 0.245 \\ 0.88 & 0.12 \end{bmatrix}$$

$$V^3 = \begin{bmatrix} 6.5 & 20 \end{bmatrix}$$

$$H^3 = \begin{bmatrix} 0.5 & 0.358 \\ 2 & 0 \end{bmatrix}$$

$t=4$:

$$\rightarrow s_t = B, a_t = a^2, \gamma_{t+1} = 20, s_{t+1} = A$$

$$\rightarrow \pi(a^1 | A) = e^{0.5} / (e^{0.5} + e^{0.358})$$

$$\boxed{\pi(a^1 | A) = 0.5354}$$

$$\Rightarrow \pi(a^2 | A) = e^{0.358} / (e^{0.5} + e^{0.358})$$

$$\boxed{\pi(a^2 | A) = 0.4646}$$

$\pi(\cdot | B)$ unchanged.

Date: _____

$$\rightarrow \delta_t = 20 + 0.9 \times V(A) - V(B)$$

$$= 20 + 0.9 \times 6.5 - 20 = 5.85$$

$$\rightarrow V(B) = 20 + 0.5 \times 5.85$$

$$V(B) = 22.925$$

$$\rightarrow H(B, a^2) = 0 + 0.1 \times 5.85 \times (1 - \pi(a^2 | B))$$

$$= 0.1 \times 5.85 \times (1 - 0.12)$$

$$H(B, a^2) = 0.51 \quad 0.515$$

So,

$$\pi^4 = \begin{bmatrix} 0.5354 & 0.4646 \\ 0.88 & 0.12 \end{bmatrix}$$

$$V^4 = \begin{bmatrix} 6.5 & 22.925 \end{bmatrix}$$

$$H^4 = \begin{bmatrix} 0.5 & 0.358 \\ 2 & 0.515 \end{bmatrix}$$

$t=5$:

$$\rightarrow S_t = A, a_t = a^1, \lambda_{t+1} = 10, S_{t+1} = A$$

$$\rightarrow \pi(a^1 | B) = e^2 / (e^2 + e^{0.515})$$

$$\pi(a^1 | B) = 0.815$$

$$\pi(a^2 | B) = e^{0.515} / (e^2 + e^{0.515})$$

$$\pi(a^2 | B) = 0.185$$

$\pi(\cdot | A)$ unchanged.

$$\rightarrow \delta_t = 10 + 0.9 \times V(A) - V(A)$$

$$= 10 + 0.9 \times 6.5 - 6.5 = 9.35$$

$$\rightarrow V(A) = V(A) + 0.5 \times \delta_t = 6.5 + 0.5 \times 9.35$$

$$V(A) = 11.175$$

$$\rightarrow H(A, a^1) = 0.5 + 0.1 \times 9.35 \times (1 - \pi(a^1 | A))$$

$$H(A, a^1) = 0.9344$$

So,

$$\pi^S = \begin{bmatrix} 0.5354 & 0.4646 \\ 0.815 & 0.185 \end{bmatrix}$$

$$v^S = \begin{bmatrix} 11.175 & 22.925 \end{bmatrix}$$

$$M^S = \begin{bmatrix} 0.9344 & 0.358 \\ 2 & 0.515 \end{bmatrix}$$

The final policy will be,

$$\pi(a'|A) = e^{0.9344} / (e^{0.9344} + e^{0.358}) \\ = 0.64$$

$$\pi(a^2|A) = 0.36$$

$\pi(\cdot | B)$ unchanged

Thus,

$$\pi^B = A \begin{bmatrix} a' & a^2 \\ 0.64 & 0.36 \\ 0.815 & 0.185 \end{bmatrix}$$

x

x

x